

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

APPLICATION FOR LETTERS PATENT

09/696707

**System and Method for Designing a Logical Model of a  
Distributed Computer System and Deploying Physical  
Resources According to the Logical Model**

Inventor(s):

Robert Welland

Galen Hunt

Aamer Hydrie

Steven Levi

Jakob Rehof

Bassam Tabbara

Mark VanAntwerp

1 **TECHNICAL FIELD**

2 This invention relates to distributed computer systems, such as Internet-  
3 based Services or Websites. More particularly, this invention pertains to a way to  
4 convert a logical, scale-independent model of an application for a distributed  
5 computer system to an actual physical configuration.  
6

7 **BACKGROUND**

8 It is no secret that Internet usage has exploded over the past few years and  
9 continues to grow rapidly. People have become very comfortable with many  
10 services offered on the World Wide Web (or simply "Web"), such as electronic  
11 mail, online shopping, gathering news and information, listening to music,  
12 viewing video clips, looking for jobs, and so forth. To keep pace with the growing  
13 demand for Internet-based services, there has been tremendous growth in the  
14 computer systems dedicated to hosting Websites, providing backend services for  
15 those sites, and storing data associated with the sites.

16 One type of distributed computer system is an Internet data center (IDC),  
17 which is a specifically designed complex that houses many computers for hosting  
18 Internet-based services. IDCs, which also go by the names "Webfarms" and  
19 "server farms", typically house hundreds to thousands of computers in climate-  
20 controlled, physically secure buildings. These computers are interconnected to run  
21 one or more programs supporting one or more Internet services or Websites. IDCs  
22 provide reliable Internet access, reliable power supplies, and a secure operating  
23 environment.

24 Fig. 1 shows an Internet data center 100. It has many server computers 102  
25 arranged in a specially constructed room. The computers are general-purpose

1 computers, typically configured as servers. An Internet data center may be  
2 constructed to house a single site for a single entity (e.g., a data center for Yahoo!  
3 or MSN), or to accommodate multiple sites for multiple entities (e.g., an Exodus  
4 center that host sites for multiple companies).

5 The IDC 100 is illustrated with three entities that share the computer  
6 resources: entity A, entity B, and entity C. These entities represent various  
7 companies that want a presence on the Web. The IDC 100 has a pool of additional  
8 computers 104 that may be used by the entities at times of heavy traffic. For  
9 example, an entity engaged in online retailing may experience significantly more  
10 demand during the Christmas season. The additional computers give the IDC  
11 flexibility to meet this demand.

12 While there are often many computers, an Internet service or Website may  
13 only run a few programs. For instance, one Website may have 2000-3000  
14 computers that run only 10-20 distinct software components. Computers may be  
15 added daily to provide scalability as the Website receives increasingly more  
16 visitors, but the underlying programs change less frequently. Rather, there are  
17 simply more computers running the same software in parallel to accommodate the  
18 increased volume of visitors.

19 Today, there is no conventional way to architect Internet Services in a way  
20 that abstracts the functionality of the Service from the underlying physical  
21 deployment. Little thought has gone into how to describe a complete Internet  
22 Service in any manner, let alone a scale-invariant manner. At best, Internet  
23 Service operators might draft a document that essentially shows each and every  
24 computer, software program, storage device, communication link, and operational  
25 relationship in the Website as of a specific time and date. The downside with such

1 physical schematics is, of course, that the document is always out of date, must be  
2 updated as the Service grows in physical resources and hence, it is of limited  
3 usefulness as a management tool. Furthermore, while a human may understand  
4 such a document, it holds no meaning to a computer.

5 Moreover, managing the physical resources of the distributed computer  
6 system for an Internet Service is difficult today. Decisions such as when to add (or  
7 remove) computers to carry out portions of the Internet Service are made by  
8 human operators. Often times, these decisions are made based on the operators'  
9 experience in running the Internet Service. Unfortunately, with the rapid growth  
10 of services, there is a shortage of qualified operators who can make real-time  
11 decisions affecting the operation of a Website. Accordingly, it would be beneficial  
12 if some of the managerial aspects of running a Internet service could be  
13 automated.

## 14 15 SUMMARY

16 A system facilitates the design and implementation of large-scale  
17 distributed computer applications, such as Internet Services and Websites. The  
18 applications are implemented as software distributed over many interconnected  
19 computer nodes, such as server data centers, Internet data centers (IDCs), Web  
20 farms, and the like.

21 The system has a modeling system and a deployment system. The  
22 modeling system permits developers to architect the hardware and software used  
23 to implement the applications in an abstract manner. The modeling system defines  
24 a set of components used to describe the functionality of an application in a  
25 logical, scale-independent manner. In the described implementation, the modeling

1 system defines several model components: a module, a port, and a wire. The  
2 model also admits an unlimited set of model extensions including, but not limited  
3 to stores, event sources, event sinks, and event wires.

4 The module is the basic functional unit and represents a container of  
5 behavior that may be implemented by one or more computers running one or more  
6 software programs. For instance, in the context of a Website, one module might  
7 represent a front end that renders HTML pages, another module might represent a  
8 login database, and another module might represent a mailbox program. A port is  
9 a service access point for the module. All communications into and out of the  
10 module goes through a port. A wire is the logical binding that defines an allowed  
11 communication route between two ports.

12 While the model consists of the three basic components described above  
13 (namely modules, ports, and wires), the model can be augmented with numerous  
14 extensions, specializations of the basic components. For example, a store is a  
15 basic unit of storage and a specialization of the module. A store represents a  
16 logical amount of storage, which may be implemented by any number of physical  
17 disks or other storage media. Like the module, the store represents behavior, in  
18 the case the ability to save and retrieve data. Also like the module, the store can  
19 communicate with other modules and stores through ports and wires. A store  
20 differs from a module in that it is labeled with additional attributes such as the  
21 amount of storage required, required access speed, or a minimum number of  
22 outstanding storage requests. The store extends the model by adding a specialized  
23 type of module with additional semantic information.

24 The model can be further augmented with ports extensions. For example,  
25 an event source and an event sink are used for discrete semantic messaging

1 between modules and module extensions, such as stores. Event sinks are  
2 specialized ports in that they are communication access points between model  
3 components, but with additional semantics, namely the specific events.

4 The model can also be augmented with wires extensions. For example, an  
5 event wire is a logical connection between event sources and event sinks, and  
6 carries event messages used to inform modules and implement policy. While most  
7 wire extensions allow communication at run time, it is possible for some wire  
8 extensions to transfer data only at compile or initialization time.

9 The model components are arranged and interconnected to form a scale-  
10 independent model of the application. Each component specifies some  
11 functionality of the application.

12 Once a logical model is created, the deployment system uses the logical  
13 model to automatically deploy various computer/software resources to implement  
14 the application. The deployment system converts each of the model components  
15 into one or more instances that correspond to physical resources. As one example,  
16 the resources correspond to computer nodes of a distributed computer system that  
17 are loaded with specific types of software to implement the function represented  
18 by the model components. The deployment system initially installs the application  
19 and then dynamically and automatically modifies the resources used to implement  
20 the application in an ongoing basis as the operating parameters of the application  
21 change.

22 In one implementation, the deployment system includes a service running  
23 state to store the logical model and track instances of the model components as  
24 they are created (or destroyed). A resource manager tracks the computer nodes  
25 available for allocation and tracks the nodes as they are allocated to correlate the

nodes with the instances. The deployment system further includes a loader to load software onto newly allocated computer nodes to implement the logical functions represented by the model components.

#### **BRIEF DESCRIPTION OF THE DRAWINGS**

Fig. 1 illustrates a conventional Internet data center (IDC).

Fig. 2 illustrates a set of model components that form the building blocks for modeling an Internet Service, along with the associated schema.

Fig. 3 illustrates a database application for an Internet Service that is modeled in terms of the components.

Fig. 4 illustrates an Internet-based email Internet Service.

Fig. 5 is a block diagram of a computer that may be used to implement the modeling software for modeling an Internet Service.

Fig. 6 is a flow diagram of a process for modeling an Internet Service.

Fig. 7 is a block diagram of a deployment system that converts a logical model to a fully functioning physical implementation

Fig. 8 illustrates a translation of the logical model into real-world instances.

Fig. 9 illustrates exemplary data records of an instance database used to store the real-world instances.

Fig. 10 is a flow diagram of a process for deploying resources for the application based on the logical model.

## DETAILED DESCRIPTION

A design system for designing applications of distributed computer systems includes a modeling system and a deployment system. The modeling system permits developers of applications for distributed computer systems (e.g., server data centers, Internet data centers (IDCs), Web farms, and the like) to architect the hardware and software in an abstract manner. The modeling system defines a set of components used to describe the functionality of an application in a logical, scale-independent manner. An "application" within this context refers to an entire service hosted on the distributed computers. For instance, an Internet data center may host a Website for an online retailer, where the application entails the entire software and hardware configuration that implements the online retailer's Internet presence. The application might include, for example, a front end to handle client requests, an order processing system, a billing system, an inventory system, and a database system.

The model components are arranged and interconnected to form a scale-independent model of the application. Each component specifies some functionality of the application. The model can then be used to construct a scalable physical blueprint in terms of which machines run which pieces of software to form the application.

The deployment system uses the logical model to deploy various computer/software resources in real-time as the applications need them. The deployment system converts each of the model components into one or more instances that correspond to physical resources. The deployment system tracks the instances and all available resources. The deployment system decides when resources should be added (or removed) and monitors the current state of



1 implementation. The deployment system installs the application and then  
2 dynamically and automatically modifies the resources used to implement the  
3 application in an ongoing basis as the operating parameters of the application  
4 change.

5 The design system is described in the context of Internet Services and  
6 Websites, such as might be deployed in Internet data centers, because modeling  
7 Internet Services represents one suitable use of the system. However, the design  
8 system may be implemented to model other large size and scalable applications for  
9 computer systems. Accordingly, the design system can be implemented in a wide  
10 variety of ways, including Internet-based implementations and non-Internet-based  
11 implementations.

### 12 13 **Model Components and Schema**

14 The modeling system defines several model components that form the  
15 building blocks of a logical, scale-independent application: a module, a port, and a  
16 wire. It also defines a set of model extensions including, but not limited to: a  
17 store, an event source, an event sink, and an event wire. In a design tool, the  
18 components are represented pictorially as graphical elements or symbols that may  
19 be arranged and interconnected to create scale-independent models of Website  
20 applications. The graphical elements have an associated schema that dictates how  
21 the functional operations being represented by the graphical elements are to be  
22 specified.

23 Fig. 2 illustrates a set of model components 200 that form the building  
24 blocks of logical, scale-independent Internet Services. The components include a  
25 module, as represented by modules 202(A)-202(C), ports 206, wires 208, and

1 extensions such as a store 204, event sources 210, event sinks 212, and event wires  
2 214. The components 200 are arranged in a no particular manner other than to  
3 foster discussion of their individual traits.

4 A module 202 represents a basic unit of functionality for the Internet  
5 Service. It is a logical entity that represents some portion of the application as  
6 might be deployed at the IDC, but it does not necessarily have a physical  
7 manifestation. The module often corresponds to a software program that handles a  
8 logical set of tasks for the Service. For instance, one module might represent a  
9 front end for a Website, another module might represent a login database, and  
10 another module might represent an electronic mail program.

11 Each module 202 is a container of behavior. A simple module is indivisible  
12 and has associated a unique identifier. Modules can be nested into a hierarchy of  
13 modules to form more complex behaviors. In a module hierarchy, the leaf  
14 modules are simple modules, and the non-leaf modules are compound modules.

15 Each module 202 defines a unit of scaling. While one module logically  
16 represents a functional operation of the Service, the module may be deployed to  
17 any number of computers when actually implemented. In this way, the module is  
18 scale-independent, allowing the number of underlying computers used to  
19 implement the module to change at over time. When converted to a physical  
20 implementation, "module instances" are created from the modules. The module  
21 instances are assigned a unique identifier and maintain ancestral data regarding  
22 which module created them. The module instances of simple modules are often  
23 called "engines", which are software programs that run on individual computer.

24 Extensions to the model are additional components that specialize the role,  
25 behavior, and possibly graphical representation of the base components.

1 Exemplary extensions include, but are not limited to, store 204, event source 210,  
2 event sink 212, and event wire 214.

3 A store 204 is the most basic unit of storage. It represents a logical storage  
4 partition, which may be implemented by any number of physical disks or other  
5 storage media.

6 A port 206 is a service access point (SAP) for a module 202 or store 204.  
7 All service-related communications into and out of a module go through a port  
8 206. Each port 206 has a "type", which is a set of attributes describing format,  
9 semantics, protocol, and so forth. At runtime, the port represents a set of physical  
10 ports associated with the instantiated engines of the modules. Note that a given  
11 module might have any number of ports representing different services or  
12 functionality provided by the module.

13 A wire 208 is the logical binding that defines an allowable communication  
14 route between two ports 206. Each wire 208 can be type-checked (i.e., with  
15 respect to protocols, roles) and defines protocol configuration constraints (e.g.,  
16 HTTP requires TCP, TCP requires IP, etc.).

17 Event sources 210 and event sinks 212 are used for discrete semantic  
18 messaging between modules and module extensions, such as stores. An event  
19 wire 214 is a logical connection between sources and sinks, and carries event  
20 messages used to inform modules or module extensions and implement policy  
21 (e.g., scaling, fail-over, monitoring, application processes, etc.).

22 The event sources 210 and event sinks 212, together with the ports 206,  
23 collectively form interfaces for communications to and from the modules 202 and  
24 module extensions, such as stores 204. The event sources and sinks may be  
25 implemented as ports that are configured for message handling.

1 The model components 200 are depicted as graphical icons or symbols that  
2 may be selected and interconnected using a modeling system (described below in  
3 more detail). In the illustrated example, the modules 202 are depicted as blocks,  
4 the store 204 is depicted as a disk storage icon, and the ports 206 are depicted as  
5 spherical knobs projecting from the modules or module extensions, such as stores.  
6 Additionally, the wires 208 are depicted as bold lines, the event sources 210 are  
7 depicted as triangles pointing away from the module or module extension, the  
8 event sinks 212 are depicted as triangles pointing toward the module or module  
9 extension, and the event wire 214 is depicted as a dashed line.

10 The graphical icons have an associated schema that dictates how the  
11 functional operations being represented by the icons are to be specified. For  
12 instance, a module icon may have a predefined schema that specifies the hardware  
13 and software resources used to implement the functionality represented by the  
14 module. Thus, a module for a database function might have characteristics  
15 pertaining to the kind of database (e.g., relational), the data structure (e.g., tables,  
16 relationships), software (e.g., SQL), software version, and so forth.

17 Fig. 2 also illustrates the schema underlying the graphical elements as  
18 exemplary data structures associated with the model components. Module 202(A)  
19 has an associated structure 220 that contains various characteristics for the  
20 module, such as functionality, processing requirements, software, and so forth.  
21 Modules 202(B) and 202(C) have similar structures (not shown). Model  
22 extensions also have associated structures. The store 204 has a corresponding  
23 structure 222 that defines the requirements for storage. The store schema structure  
24 222 might include, for example, the kind of storage (e.g., disk), the storage format,  
25 and so on.

09556707 " 403400

Each port 206 has a schema structure, as represented by structure 224, which dictates the port's type. Each wire 208 is also associated with a schema structure, such as structure 226, which outlines the protocols implemented by the connection. Similar schema structures may also be provide for event sources event sinks, and event wires.

Using the model components, a developer can logically describe and configure scale-independent Internet Service prior to physically laying them out in Internet data centers. The developer drafts a model using a user interface to select and interconnect the model components. Once constructed, the modeling software generates the Internet Service based on the depicted model and the underlying schema. The Service may subsequently be converted into a physical blueprint that details the computers and software needed to implement the Service for a specified number of clients.

The scale-invariant nature of the modeling system allows Internet Service developers to focus only on designing software for a specific functional task (e.g., front end, login database, email program, etc.). All external communications can then be expressed in terms of transmitting to and receiving from one or more associated ports. In this manner, the Service developers need not worry about how many machines will be used to run the module, or how other modules of the scale-independent Internet Service are being configured.

### **Exemplary Module and Application**

Fig. 3 shows a fault-tolerant SQL (structure query language) database module 300 to demonstrate how the model components may be organized and connected to represent a portion of an application. In this example, the database

1 module 300 represents a SQL database that may be used independently or as a  
2 component in a larger application. The SQL database module 300 has a module  
3 interface composed of a single port 302 that implements the TDS (Tabular Data  
4 Stream) protocol.

5 The SQL database module 300 is a compound module made up of three  
6 simple modules: a fail-over policy module 310, a primary SQL module 312, and a  
7 secondary SQL module 314. The primary and secondary SQL modules represent  
8 dual programs that operate in parallel so that, in the event that the primary module  
9 312 crashes, the secondary module 314 can assume the role without loss of  
10 service. The database module 300 also has a data store 316 that represents the  
11 memory storage for the SQL database module.

12 The primary SQL module 312 has a module interface that includes a first  
13 port 320 for communicating with the compound module port 302 and a second  
14 port 322 for communicating with the store 316. The primary SQL module 312  
15 also has an event source 324 and an event sink 326 for handling event messages  
16 from the fail-over policy module 310. Similarly, the secondary SQL module 314  
17 has a module interface with a first port 330 for communicating with the compound  
18 module port 302, a second port 332 for communicating with the store 316, and an  
19 event sink 334 for receiving events from the fail-over policy module 310. A wire  
20 336 interconnects the external compound module port 302 with the ports 320 and  
21 330 of the primary and secondary SQL modules, respectively.

22 The store 316 has a port 340 to communicate with the primary and  
23 secondary SQL modules 312 and an event sink 342 to receive event messages  
24 from the fail-over policy module 310. A wire 344 interconnects the store port 340  
25

1 with the ports 322 and 332 of the primary and secondary SQL modules,  
2 respectively.

3 The fail-over policy module 310 has a module interface that includes three  
4 event sources and one event sink. An event sink 350 receives a "fail" event from  
5 the event source 324 of the primary SQL module 312 via an event wire 352 when  
6 the primary SQL module experiences some failure. In response to receiving a  
7 "fail" event, the fail-over policy module 310 concurrently issues a first event to  
8 stop the failed primary module 312, another event to assign the secondary module  
9 314 as the new owner of the store 316, and a third event to start the secondary  
10 module 314. The "stop" event is issued via an event source 354 over an event  
11 wire 356 to the event sink 326 of the primary SQL module 312. The "stop" event  
12 directs the primary SQL module 312 to halt operation.

13 The fail-over policy module 310 issues an "assign owner" (AO) event from  
14 an event source 358, over the event wire 360 to the event sink 342 of the store  
15 316. The assign owner event directs the storage mechanisms to switch to allowing  
16 access by the secondary SQL module 314, rather than the primary SQL module  
17 312. The fail-over policy module 310 also issues a "start" event from event source  
18 362 over event wire 364 to the event sink 334 of the secondary module 314. The  
19 start event directs the secondary SQL module to start operation in place of the  
20 primary SQL module.

21 The SQL database module 300 illustrates how the base model components  
22 and exemplary model extensions—modules, ports, wires, stores, event sources,  
23 event sinks, and event wires—may be arranged and interconnected to form a  
24 complex module. The developer specifies the characteristics associated with each  
25 component according to the prescribed schema. The complex module may in turn

1 be added to other simple or complex modules to form other complex modules.  
2 Eventually, the largest complex module becomes the Internet Service, which may  
3 then be used to form a blueprint for deploying to the data center.

4 Fig. 4 shows a simplified application 400 for an online retailer. The  
5 application 400 includes a front end module 402, a catalog module 404, an order  
6 processing module 406, and a fulfillment module 408. The application 400 also  
7 includes a customer database 410 and the fault-tolerant SQL database module 300.  
8 Notice that the SQL database module 300 is the same as that shown in Fig. 3 to  
9 illustrate how complex modules can be nested into even greater complex modules  
10 to form an application.

11 The front end module 402 handles requests from clients who wish to shop  
12 with the online retailer. The front end module 402 has a port 420 that  
13 accommodates communications with external clients using the TCP/IP protocol  
14 over the Internet. The front end module 402 also has an order port 422 to define a  
15 communication exchange with the order processing module 406 and a catalog port  
16 424 for communication flow to the catalog module 404. The ports 422 and 424  
17 may be configured according to any of a variety of types, which support any of a  
18 number of protocols including SOAP, TCP, or UDP. An event sink 426 is also  
19 provided to receive a "new product" message from the catalog module 404 when a  
20 new product has been added to the catalog.

21 The catalog module 404 provides catalog information that may be served  
22 by the front end to the requesting clients. The catalog module 404 has a front end  
23 port 430 connected via a wire 432 to the catalog port 424 of the front end module  
24 402. The front end port 430 has a type that matches the catalog port 424. The  
25 catalog module 404 also has an event source 434 for communicating the "new



product" messages over wire 436 to the event sink 426 of the front end module 402.

A SQL port 438 interfaces the catalog module 404 with the SQL database module 300. The SQL port 438 has a type that utilizes the TDS protocol for the communication exchange with the external port 302 of the SQL database 300.

The order processing module 406 has a front end port 440 to define a communication interface with the front end module 402 via a wire 442. The order processing module 406 also has a fulfillment port 444 to facilitate communication with the fulfillment module 408 over wire 446 and a database port 448 to facilitate communication with the customer database 410 via wire 450.

An event source 452 is provided at the order processing module 406 to pass "order complete" events to the fulfillment module 408 via wire 454. These events inform the fulfillment module 408 that an order is complete and ready to be filled. A second event source 456 passes "new account" events to the customer database 410 via wire 458 whenever a new customer orders a product.

The fulfillment module 408 has an order port 460 to provide access to the wire 446 to the order processing module 406 and a database port 462 to interface with the customer database 410. The fulfillment module 408 also has an event sink 464 to receive the "order complete" events from the order processing module 406.

The customer database 410 has an order port 470 to provide access to wire 450 and a fulfillment port 472 to facilitate communication with the fulfillment module 408 via wire 474. The customer database 410 further has an event sink 476 to receive the "new account" events from the order processing module 406.

1 The modeling approach illustrated in Figs 3 and 4 is tremendously  
2 beneficial because it allows developers and IDC operators to view the entire  
3 Internet Service in terms of functional pieces independent of deployment scale.  
4 The online retailer Internet Service 400, for example, requires a front end unit, a  
5 catalog unit, an order processing unit, and a fulfillment unit regardless of whether  
6 the retailer is handling 100 hits a day or 10 million hits per day.

7 The scale-independent nature frees the developer to focus on his/her little  
8 piece of the Service. For instance, a developer assigned the task of building the  
9 front end module 402 need only be concerned with writing software code to  
10 facilitate response/reply exchanges. Any communication to and from the module  
11 is defined in terms of order-related data being passed to the order processing  
12 module 406 via the order port 422 and product data being received from the  
13 catalog module 404 via the catalog port 424. The developer defines the data flow  
14 to and from the order port 422 and the catalog port 424 according to their  
15 respective associated protocol types.

16 The Internet Service 400 can then be used to construct a computer system  
17 that hosts the online retailer. Initially, the online retailer may not receive very  
18 much traffic, especially if launched away from the Christmas season. So, perhaps  
19 the front end module 402 deploys initially to only a few computers to handle the  
20 light traffic from the Internet. But, suppose that over time the site becomes more  
21 popular and the Christmas season is fast approaching. In this situation, the online  
22 retailer may authorize the IDC operator to add many more computers for the front  
23 end tasks. These computers are equipped with software and configured to accept  
24 HTTP requests for product information and to serve Web pages containing the  
25

1 product information. The computers are added (or removed) as needed, without  
2 altering the basic description of the Internet Service 400.

### 3 4 **Computer-Based Modeling System and Method**

5 Fig. 5 shows an exemplary computer system 500 that implements modeling  
6 software used to design Internet Services. The modeling computer may be  
7 implemented as one of the nodes in a Internet Service, or as a separate computer  
8 not included as one of the nodes. The modeling computer has a processor 502,  
9 volatile memory 504 (e.g., RAM), and non-volatile memory 506 (e.g., ROM,  
10 Flash, hard disk, optical, RAID memory, etc.). The modeling computer 500 runs  
11 an operating system 510 and modeling system 512.

12 For purposes of illustration, operating system 510 and modeling system 512  
13 are illustrated as discrete blocks stored in the non-volatile memory 506, although it  
14 is recognized that such programs and components reside at various times in  
15 different storage components of the computer 500 and are executed by the  
16 processor 502. Generally, these software components are stored in non-volatile  
17 memory 506 and from there, are loaded at least partially into the volatile main  
18 memory 504 for execution on the processor 502.

19 The modeling system 512 includes a user interface 514 (e.g., a graphical  
20 UI) that presents the pictorial icons of the model components 516 (e.g., modules,  
21 ports, sources, sinks, etc.), a component schema database 518, a logical-to-  
22 physical converter 520, and an instance-tracking database 522. The modeling  
23 system 512 allows a developer to design an Internet Service by defining modules,  
24 ports, wires, and event message schemes. The user interface 514 presents symbols  
25 of the components 516, such as the symbols shown in Figs 2-4, and permits the

1 developer to arrange and interconnect them. The UI 514 may even support  
2 conventional UI techniques as drag-and-drop operations.

3 The symbols depicted on the screen represent an underlying schema 518  
4 that is used to define the model. For instance, a block-like module symbol is  
5 associated with the characteristics of the functionality that the module is to  
6 represent in the Internet Service. Thus, the developer may define a database  
7 module that has characteristics pertaining to the kind of database (e.g., relational),  
8 the data structure (e.g., tables, relationships), software (e.g., SQL), software  
9 version, and so forth. Accordingly, by drafting the model on the UI, the developer  
10 is architecting the entire schema that will be used to design the scale-independent  
11 Internet Service.

12 Once the Internet Service is created, the logical-to-physical converter 520  
13 converts the Service to a physical blueprint that details the number of computers,  
14 software components, physical ports, and so forth. The converter takes various  
15 parameters—such as how many site visitors are expected, memory requirements,  
16 bandwidth requirements, processing capabilities, and the like—and scales the  
17 Internet Service according to the schema 518 created by the developer. The  
18 converter 520 specifies the number of computers needed to implement each  
19 module, the number of disks to accommodate the stores, and the types of  
20 communications protocols among the modules and stores. The identity of every  
21 component instance is recorded in an instance-tracking database 522. Instances in  
22 the instance-tracking database 522 include those for modules, port, wires, and  
23 instances of model extensions such as stores, event ports, and event wires.

24 In one embodiment, the developer writes management policy, which issues  
25 commands on the schema 518 to create new instances of modules, port, and wires

1 to deploy the Internet Service. Developers may choose to write management  
2 policy instead of using fully automatic logical-to-physical converter 520 when  
3 they want finer control over the growth and management of the Internet Service.  
4 The management code issues commands using the namespace defined by the  
5 schema 518, but the commands operate on individual module, port, and wire  
6 instances. The commands are still dispatched through the converter 520, which  
7 allocates nodes to the management policy. Whether operating automatically or  
8 driven by management policy code, the convert 520 records in the instance-  
9 tracking database 522 the individual instances of modules, port, and wires.

10 In this manner, the modeling system changes the development effort from a  
11 node-centric approach for architecting Internet Services to an application-centric  
12 approach. Within conventional node-centric methodology, the focus was on the  
13 computers and how they were laid out. The Internet Service was then loaded onto  
14 the nodes in an ad hoc manner. With the new application-centric approach, the  
15 focus is initially on the Internet Service itself. The physical nodes used to  
16 implement the Internet Service are derived in terms of the Service schema once it  
17 is specified. The instance-tracking database 522 gives both developers and  
18 operators information about how many instances of each module are running at  
19 any time and how the modules are connected using port instances and wires  
20 instances within the Service schema.

21 Fig. 6 shows a method for modeling a scale-independent Internet Service.  
22 The method 600 may be implemented, for example, by the modeling system 512  
23 executing on the modeling computer 500. In such an implementation, the method  
24 is implemented in software that, when executed on computer 500, performs the  
25 operations illustrated as blocks in Fig. 6.

At block 602, the modeling system 512 allows the developer to define the modules and extensions, such as stores, that form the functional elements of the Internet Service. The UI 514 enables the developer to create modules and extensions, such as stores, and to define their characteristics as prescribed by a predetermined schema. This entry process begins to construct the logical building blocks of the Service.

At block 604, the modeling system 512 enables the developer to define the ports for the modules and module extensions, such as stores. The developer selects the type of ports. The modeling system ensures compatibility of ports that are connected to one another. At block 606, the developer also defines other extensions, such as events that may be passed among modules and module extensions, such as stores. For example, the developer creates event sources and event sinks to accommodate the various events. At block 608, the developer uses the modeling system 512 to interconnect the ports with wires and the port extensions, such as event sources/sinks with wire extensions, such as event wires. By joining the various modules and module extensions, such as stores, the developer effectively forms a logical representation of the Internet Service.

At block 610, the modeling system 512 generates an Internet Service using the graphical representation constructed by the developer through its associated schema. The modeling system 512 generates the logical specifications associated with the graphical model, including the characteristics of the modules and module extensions, such as stores, as well as the types of the ports and port extensions, such as event sources/sinks. The Internet Service provides a complete logical representation of the Service that will eventually be implemented at the Internet

1 data center. The Internet Service description may be stored on disk or some other  
2 form of computer-readable medium (block 612).

3 At block 614, the modeling system 512 converts the Internet Service  
4 description to a physical blueprint that specifies the computers, the software run  
5 by each of the computers, and the interconnections among the computers. This  
6 physical blueprint may be used by the operator to install and manage the Internet  
7 Service.

### 8 9 **Computer-Based Deployment System and Method**

10 Once a logical model is created, an automatic computer-based deployment  
11 system uses the logical model to deploy various computer/software resources to  
12 implement the Internet Service. The deployment system converts each of the  
13 model components into one or more instances that correspond to physical  
14 resources, such as nodes of a distributed computer system that are loaded with  
15 specific types of software to implement the function represented by the model  
16 components. The deployment system initially installs an Internet Service on the  
17 physical resources according to the logical model. It then dynamically and  
18 automatically modifies the resources used to implement the Internet Service in an  
19 ongoing basis as the operating parameters of the application change.

20 In one embodiment, the deployment system installs an Internet Service  
21 under the direction of management policy code authored by the Service  
22 developers. When operated under the direction of management code, the  
23 deployment system creates instances as instructed by the management code. The  
24 management code assumes partial or complete responsibility for monitoring the  
25 Service and determining when instances should be created and destroyed.

Fig. 7 shows a deployment system 700 that converts the logical model to a fully functioning physical implementation. The deployment system 700 includes management policy 702, a core logical-to-physical converter 704, and hardware/software resources 706 that are all interconnected via a wireless and/or wire-based communication network 708 (e.g., a LAN, a WAN, intranet, Internet, combinations thereof, etc.). In this example, the hardware/software resources are illustrated as computer nodes of a distributed computer system, as represented by computers 706(1), 706(2), ..., 706(N). The management policy 702 and core converter 704 may be implemented on one or more computers, which may or may not be part of the nodes in the distributed computer system.

For purposes of discussion, the deployment system 700 is described in the context of an Internet Service that is executed at an Internet data center having an abundance of generic computer nodes. The nodes can be allocated to one or more Internet Services from a reserve pool of nodes, as illustrated in Fig. 1.

The management policy 702 implements one or more policies devised by the developer or operator of the Internet Service. The policies specify when instances derived from the logical model should be created, manipulated, and destroyed. The management policy monitors various events generated by the nodes and implements policy decisions regarding how to handle the events. By specifying when and what instances should be created (or destroyed), the management policy 702 effectively dictates when hardware/software resources 706 should be added (or removed) to support the changing demands of the Internet Service.

The core converter 704 implements the policy decisions made by the management policy 702. The runtime converter 704 has a service running state



710 that tracks all instances of the model components currently in existence. That is, the service running state 710 tracks the elements in the physical world with respect to the logical model. The service running state 710 maintains a copy of the logical model, such as online retailing model 400 of Fig. 4. The logical model is created by the modeling system described above with respect to Figs. 5 and 6. The current instances are maintained in the instance-tracking database 522. The records in instance-tracking database 522 include such information as identify of the instance, the name of the logical component from which it is derived, the node on which it is running, the network addresses representing the ports of the modules and module extensions, such as stores, type of software loaded on the node, various protocols supported by the instance, and so forth. The instance-tracking database 522 tracks not only module instances, but also port instance, wire instances, and can also track instances of extensions such as stores, event ports, and event wires.

The instances are derived from the logical model. The management policy 702 articulates the number of instances of each model component used to implement the Internet Service at any given time. For example, suppose the Internet Service requires one hundred computers to effectively implement a front end that handles site traffic at 99.9% efficiency with each computer running at 70% utilization. The management policy 702 might further specify that more computers should be added if some policy threshold is met (e.g., efficiency rating drops below some threshold or computer utilization rises above some threshold) or removed if another threshold is met.

A resource manager 716 tracks all of the physical resources available to the Internet Service. These resources include computer nodes, storage, software, and

so forth. Records identifying all of the resources are kept in the resource database 718. For instance, there might be one record for each computer node, storage device, and software module in the Internet data center. The records contain such information as the identity of the allocated nodes, computing characteristics and capabilities, the application(s) to which they are allocated, the date and time of allocation, and so forth.

The resource manager 716 allocates the resources as needed or requested by the Internet Service according to the policy implemented by the policy manager 702. The allocation depends upon the availability of resources at the time of request. The resource manager 716 may also recover resources that are no longer needed by the Internet Service and return the resources to the pool of available resources.

Upon allocation (or recovery) of a resource, the resource manager 716 posts a record to the resource database 718 reflecting which resource is allocated to (or recovered from) which Internet Service. As an example, when an Internet Service desires more nodes for the front end tasks, the resource manager 716 allocates one or more free nodes from the pool of resources to the Internet service.

The core logical-to-physical converter 704, manages the creation of new instances in the physical world from the model components specified in the logical model 400 through a loader 722. A loader 722 carries out the configuration of newly allocated resources to the functions dictated by the new instances. In this manner, when another instance of a component is desired, the management policy 720 communicates with the resource manager 716 to allocate a node from the resource pool and with the loader 722 to load the appropriate software programs onto the node.

1 The node loader 732 performs the actual loading tasks specified by the  
2 loader 722 in the core converter 704. It is the local code on an otherwise generic  
3 node to which the runtime loader 722 may communicate when configuring the  
4 node.

5 Fig. 8 shows an example of a portion of logical model 400 being converted  
6 into actual instances. To illustrate the conversion, the front end module 402 and  
7 the order processing module 406 are extracted from the online retailer service  
8 application 400 (Fig. 4). A wire 442 interconnects the two modules by logically  
9 coupling the ports 422 and 440.

10 The front end module 402 in the logical model translates to one or more  
11 computer nodes in the physical world that runs software for handling client  
12 queries. These physical instances are represented diagrammatically by the  
13 rectangles with ovals. Here, based on a policy, the front end module 402 is  
14 converted into multiple front end instances 800(1), 800(2), 800(3), ..., 800(J),  
15 which each corresponds in one-to-one fashion with a computer node loaded with  
16 software used to implement the front end of the service. The order processing  
17 module 406 translates to one or more computer nodes that run a program for  
18 processing client orders. In Fig. 8, the order processing module 406 is converted  
19 into plural order processing instances 802(1), 802(2), ..., 802(K).

20 The ports 422 and 440 represent network ports with protocols and roles  
21 within protocols in the logical world. They translate to the physical world as a set  
22 of network addresses used to communicate with the software at the respective  
23 modules. For instance, the physical translation of a logical port might be IP Port  
24 80, using HTTP (hypertext transport protocol). In Fig. 8, one logical port in the  
25 logical model is converted to a port address for each instance of the module.

1 Logical port 422 converts into physical address ports 804(1)-804(J) and logical  
2 port 440 converts into physical ports 806(1)-806(K).

3 The wire 442 represents a set of ports with allowable communication  
4 connections. It translates to a physical mesh of all possible communication wires  
5 between the instances of each port. The maximum number of physical lines  
6 created from one wire 442 is determined as the cross product of the number of  
7 instances of each port. In Fig. 8, the wire 442 can convert to a physical connection  
8 between every port instance 804(1)-804(J) of the front end module 402 and every  
9 port instance 806(1)-806(K) of the order processing module 406. The number of  
10 actual physical connections created from the wire 442 is determined by the  
11 management policy 702.

12 Fig. 9 illustrates exemplary records 900 in the instance-tracking database  
13 522 that tracks the instances derived from the logical model 400. In this example,  
14 the database is a relational database that stores records in tables, and the records  
15 may be linked to one another via relationships. Here, there are three tables: a  
16 module table 902, a port table 904, and a wire table 906. Each table holds one  
17 record for a corresponding instance of the logical model. Each record has a  
18 number of fields relating to the type of information that is being tracked for each  
19 instance.

20 The module table 902 tracks instances of modules in the logical model.  
21 There is one record for each module instance. Thus, with respect to the front end  
22 module of Fig. 8, there are "J" records in the module table 902 corresponding to  
23 the "J" front end instances 800(1)-800(J). Each record in the module table 902  
24 contains an instance ID to identify the individual instance, an identity of the  
25 module component from which the instance is derived, a node ID to identify the

1 computer node to which the instance is associated, the type of software loaded on  
2 the node to implement the module functionality, a software ID, an identity of the  
3 various ports to the modules, and various protocols supported by the module  
4 instance. It is noted that other implementations may include additional fields or  
5 fewer fields than those illustrated.

6 The port table 904 tracks instances of the port in the logical model, such as  
7 ports 422 and 440 in Fig. 8. A record from this table includes such information as  
8 the port ID, the model component identity, a node ID, the network address  
9 represented by the port, the instance ID of the corresponding instance, the protocol  
10 used by the port, and an ID of the wire to which the port is connected. Again,  
11 more or less fields may be used in other implementations.

12 The wire table 906 tracks instances of the wires in the logical model, such  
13 as wire 442 in Fig. 8. A record in the wire table includes a wire ID, a model  
14 component identity, the protocol supported by the wire, and information to  
15 identify each of the ports on the wire, such as a node ID, a port ID, and an instance  
16 ID.

17 Notice that the three tables can be correlated with one another via various  
18 relationships. For example, the module table 902 and the port table 904 are related  
19 by various data fields, such as the port ID field and the instance ID field. The wire  
20 table 906 correlates with the port table 904 via wire ID and port ID and with the  
21 module table 902 via instance ID. Notice also that the tables have a node ID field  
22 that provides a reference into the resource database by identifying which node the  
23 instance is associated with.

1 It is noted that the illustrated arrangement of the database is merely for  
2 discussion purposes. Many other table arrangements with more or fewer tables  
3 than illustrated may be used in other implementations.

4 Fig. 10 shows a method 1000 for deploying resources for an Internet  
5 Service based on the logical model. The method 1000 is implemented by the  
6 deployment system 700 of Fig. 7 and hence can be embodied as software that,  
7 when executed on one or more computers, performs the operations illustrated as  
8 blocks in Fig. 10.

9 At block 1002, the management policy 702 monitors various operating  
10 parameters and listens for events from the individual nodes 706 or the core  
11 logical-to-physical converter 704. The management policy 702 evaluates the  
12 parameters and events against the policy backdrop to determine whether the  
13 current physical implementation is satisfactorily supporting the Internet Service, or  
14 whether a new instance of some component should be added (block 1004). It is  
15 noted that the continuing process is described in the context of adding a new  
16 instance. However, the policy may alternatively dictate that one or more instances  
17 should be removed. Removal of instances is somewhat easier in that instances are  
18 deleted from the instance-tracking database 522 and the computer nodes are  
19 returned to the additional pool for reallocation.

20 Assuming that a new instance is desired (i.e., the "yes" branch from block  
21 1004), the management policy 702 consults the service running state 710 to  
22 understand the current number and arrangement of instances (block 1006). The  
23 management policy 702 can request various types of information of the service  
24 running state 710, such as:  
25

How many instances of a given module?  
What nodes are associated with a given model component?  
Which ports are attached to a given wire?  
What is the network address of a node?  
What wire is attached to a port on a given component?  
What components does a given component own?

Depending upon the event or operating condition, the management policy 702 can request information on a particular module in the logical model 400. For example, assume that an event has been received from a front end node indicating that the utilization has risen to above 90%. In response, a policy specifying the addition of another instance at such utilization levels is triggered. The management policy 702 asks the service running state 710 how many instances of the front end module currently exist, and information regarding how to specify an instance for the front end module.

At block 1008, the management policy 702 calls the resource manager 716 to request allocation of a new node (assuming one is available). The resource manager 716 examines the resources using data from the resource database 718 and allocates a new node that is currently available from a pool of free nodes. The resource manager 716 records the allocation in the resource database 718, identifying which node is allocated, what application it is being allocated to, the date and time that it is allocated, and so forth.

At block 1010, the management policy 702 calls the loader 722 to install software onto the allocated node and configure the node to perform the functions represented by the logical module from which the instance is created. In response, the loader 722 initializes the node by communicating with the node loader 730 of the new node via the network 708 to install the appropriate software; such an

operation might install an image of an operating system (e.g., Windows NT server operating system from Microsoft Corporation). The loader 722 then loads the appropriate software and configures the node to perform the functions represented by the logical model component from which the instance is derived. For example, for an instance of an SQL module in the logical model 400, the node loader 732 loads SQL server software. The loader 722 registers the physical port addresses with the service running state 710, which records the new instance in the instance-tracking database 522.

At block 1012, the service running state 710 is notified when the newly allocated and configured node is up and running. The service running state 710 records a new instance of the logical model in the instance-tracking database 522. The record reflects an ID of the instance, the name of the logical component from which it is derived, the allocated node, the network addresses of the node, software type and IDs, various protocols supported by the instance, and so on.

At block 1014, the management policy 702 is notified by the core runtime converter 704 that a new instance is up and running. The new instance should relieve the event or operating condition to bring operation back into compliance with the policy.

### **Conclusion**

Although the description above uses language that is specific to structural features and/or methodological acts, it is to be understood that the invention defined in the appended claims is not limited to the specific features or acts described. Rather, the specific features and acts are disclosed as exemplary forms of implementing the invention.